

ID3 Decision Tree Algorithm

Aim: To construct the Decision tree using the training data sets under supervised learning concept.

Program: Write a program to demonstrate the working of the decision tree based ID3 algorithm. Use an appropriate data set for building the decision tree and apply this knowledge to classify a new sample.

ID3 Algorithm

```
from IPython.display import Image
Image(filename='id3.png')
```

In [13]:

Out[13]:

ID3(Examples, Target_attribute, Attributes)

Examples are the training examples. Target_attribute is the attribute whose value is to be predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree. Returns a decision tree that correctly classifies the given Examples.

- Create a Root node for the tree
- If all Examples are positive, Return the single-node tree Root, with label = +
- If all Examples are negative, Return the single-node tree Root, with label = -
- If Attributes is empty, Return the single-node tree Root, with label = most common value of Target_attribute in Examples

```
Image(filename='id31.png')
```

In [2]:

Out[2]:

- Otherwise Begin
 - $A \leftarrow$ the attribute from Attributes that best* classifies Examples
 - The decision attribute for Root $\leftarrow A$
 - For each possible value, v_i , of A,
 - Add a new tree branch below Root, corresponding to the test $A = v_i$
 - Let *Examples* _{v_i} be the subset of Examples that have value v_i for A
 - If *Examples* _{v_i} is empty
 - Then below this new branch add a leaf node with label = most common value of Target_attribute in Examples
 - Else below this new branch add the subtree
 $\text{ID3}(\text{Examples}_{v_i}, \text{Target_attribute}, \text{Attributes} - \{A\})$
- End
- Return Root

* The best attribute is the one with highest information gain

```
Image(filename='entropy-infgain.png')
```

In [3]:

Out[3]:

ENTROPY:

Entropy measures the impurity of a collection of examples.

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Where, p_{+} is the proportion of positive examples in S
 p_{-} is the proportion of negative examples in S.

INFORMATION GAIN:

- **Information gain**, is the expected reduction in entropy caused by partitioning the examples according to this attribute.
- The information gain, $Gain(S, A)$ of an attribute A, relative to a collection of examples S, is defined as

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

In [4]:

```
Image(filename='trainingdata.png')
```

Out[4]:

Training Dataset:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

In [1]:

```
import pandas as pd
import numpy as np
data = pd.read_csv('PlayTennis.csv')
data
```

Out[1]:

	Day	Outlook	Temperature	Humidity	Wind	PlayTennis
0	D1	Sunny	Hot	High	Weak	No
1	D2	Sunny	Hot	High	Strong	No
2	D3	Overcast	Hot	High	Weak	Yes
3	D4	Rain	Mild	High	Weak	Yes
4	D5	Rain	Cool	Normal	Weak	Yes
5	D6	Rain	Cool	Normal	Strong	No
6	D7	Overcast	Cool	Normal	Strong	Yes
7	D8	Sunny	Mild	High	Weak	No
8	D9	Sunny	Cool	Normal	Weak	Yes
9	D10	Rain	Mild	Normal	Weak	Yes
10	D11	Sunny	Mild	Normal	Strong	Yes
11	D12	Overcast	Mild	High	Strong	Yes
12	D13	Overcast	Hot	Normal	Weak	Yes
13	D14	Rain	Mild	High	Strong	No

Image(filename='testdata.png')

In [14]:

Out[14]:

Test Dataset:

Day	Outlook	Temperature	Humidity	Wind
T1	Rain	Cool	Normal	Strong
T2	Sunny	Mild	Normal	Strong

In [3]:

```
import csv
def load_csv(filename):
    lines=csv.reader(open(filename,"r"));
    dataset = list(lines)
```

```

headers = dataset.pop(0)
return dataset,headers

```

In [4]:

```

dataset, features = load_csv('PlayTennis.csv')
dataset, features

```

Out[4]:

```

([['D1', 'Sunny', 'Hot', 'High', 'Weak', 'No'],
 ['D2', 'Sunny', 'Hot', 'High', 'Strong', 'No'],
 ['D3', 'Overcast', 'Hot', 'High', 'Weak', 'Yes'],
 ['D4', 'Rain', 'Mild', 'High', 'Weak', 'Yes'],
 ['D5', 'Rain', 'Cool', 'Normal', 'Weak', 'Yes'],
 ['D6', 'Rain', 'Cool', 'Normal', 'Strong', 'No'],
 ['D7', 'Overcast', 'Cool', 'Normal', 'Strong', 'Yes'],
 ['D8', 'Sunny', 'Mild', 'High', 'Weak', 'No'],
 ['D9', 'Sunny', 'Cool', 'Normal', 'Weak', 'Yes'],
 ['D10', 'Rain', 'Mild', 'Normal', 'Weak', 'Yes'],
 ['D11', 'Sunny', 'Mild', 'Normal', 'Strong', 'Yes'],
 ['D12', 'Overcast', 'Mild', 'High', 'Strong', 'Yes'],
 ['D13', 'Overcast', 'Hot', 'Normal', 'Weak', 'Yes'],
 ['D14', 'Rain', 'Mild', 'High', 'Strong', 'No']],
 ['Day', 'Outlook', 'Temperature', 'Humidity', 'Wind', 'PlayTennis'])

```

In [5]:

```

class Node:
    def __init__(self,attribute):
        self.attribute=attribute
        self.children=[]
        self.answer=""

```

In [6]:

```

def subtables(data,col,delete):
    dic={}
    coldata=[row[col] for row in data]
    attr=list(set(coldata))

    counts=[0]*len(attr)
    r=len(data)
    c=len(data[0])
    for x in range(len(attr)):
        for y in range(r):
            if data[y][col]==attr[x]:
                counts[x]+=1
    for x in range(len(attr)):
        dic[attr[x]]=[[0 for i in range(c)] for j in range(counts[x])]
        pos=0
        for y in range(r):
            if data[y][col]==attr[x]:
                if delete:
                    del data[y][col]
                dic[attr[x]][pos]=data[y]
                pos+=1
    return attr,dic

```

In [7]:

```

import math
def entropy(S):
    attr=list(set(S))
    if len(attr)==1:
        return 0

    counts=[0,0]
    for i in range(2):
        counts[i]=sum([1 for x in S if attr[i]==x])/(len(S)*1.0)

    sums=0

```

```

for cnt in counts:
    sums+=-1*cnt*math.log(cnt,2)
return sums

```

In [8]:

```

def compute_gain(data,col):
    attr,dic = subtables(data,col,delete=False)

    total_size=len(data)
    entropies=[0]*len(attr)
    ratio=[0]*len(attr)

    total_entropy=entropy([row[-1] for row in data])
    for x in range(len(attr)):
        ratio[x]=len(dic[attr[x]])/(total_size*1.0)
        entropies[x]=entropy([row[-1] for row in dic[attr[x]]])
        total_entropy-=ratio[x]*entropies[x]
    return total_entropy

```

In [16]:

```

def build_tree(data,features):
    lastcol=[row[-1] for row in data]
    if(len(set(lastcol))==1):
        node=Node("")
        node.answer=lastcol[0]
        return node

    n=len(data[0])-1
    gains=[0]*n
    for col in range(n):
        gains[col]=compute_gain(data,col)
    split=gains.index(max(gains))
    node=Node(features[split])
    fea = features[:split]+features[split+1:]

    attr,dic=subtables(data,split,delete=True)
    for x in range(len(attr)):
        child=build_tree(dic[attr[x]],fea)
        node.children.append((attr[x],child))
    return node

```

In [18]:

```

def print_tree(node,level):
    if node.answer!="":
        print(" "*level,node.answer)
        return

    print(" "*level,node.attribute)
    for value,n in node.children:
        print(" "*(level+1),value)
        print_tree(n,level+2)

```

In [11]:

```

def classify(node,x_test,features):
    if node.answer!="":
        print(node.answer)
        return
    pos=features.index(node.attribute)
    for value, n in node.children:
        if x_test[pos]==value:
            classify(n,x_test,features)

```

In [21]:

```

'''Main program'''
dataset,features=load_csv("PlayTennis.csv")
features = features[1:]

```

```

dataset = [ele[1:] for ele in dataset]
print("\n Features: ", features)
print("\n Dataset: ", dataset)
model=build_tree(dataset,features)

print("\n The decision tree for the dataset using ID3 algorithm is")
print_tree(model,0)
testdata,features=load_csv("testdata.csv")
for xtest in testdata:
    print("\n The test instance:",xtest)
    print("\n The label for test instance:",end="    ")
    classify(model,xtest,features)
Features:  ['Outlook', 'Temperature', 'Humidity', 'Wind', 'PlayTennis']

Dataset:  [['Sunny', 'Hot', 'High', 'Weak', 'No'], ['Sunny', 'Hot', 'High', 'Strong', 'No'], ['Overcast', 'Hot', 'High', 'Weak', 'Yes'], ['Rain', 'Mild', 'High', 'Weak', 'Yes'], ['Rain', 'Cool', 'Normal', 'Weak', 'Yes'], ['Rain', 'Cool', 'Normal', 'Strong', 'No'], ['Overcast', 'Cool', 'Normal', 'Strong', 'Yes'], ['Sunny', 'Mild', 'High', 'Weak', 'No'], ['Sunny', 'Cool', 'Normal', 'Weak', 'Yes'], ['Rain', 'Mild', 'Normal', 'Weak', 'Yes'], ['Sunny', 'Mild', 'Normal', 'Strong', 'Yes'], ['Overcast', 'Mild', 'High', 'Strong', 'Yes'], ['Overcast', 'Hot', 'Normal', 'Weak', 'Yes'], ['Rain', 'Mild', 'High', 'Strong', 'No']]

The decision tree for the dataset using ID3 algorithm is
Outlook
Rain
Wind
Strong
No
Weak
Yes
Overcast
Yes
Sunny
Humidity
Normal
Yes
High
No

The test instance: ['T1', 'Rain', 'Cool', 'Normal', 'Strong']

The label for test instance: No

The test instance: ['T2', 'Sunny', 'Mild', 'Normal', 'Strong']

The label for test instance: Yes

```

In []: